

# 3 PC  
1-10-00

IN THE U.S. PATENT AND TRADEMARK OFFICE

Applicant(s): NISHIMURA, Hideki

Application No.:

Group:

Filed: October 19, 1999

Examiner:

For: KEY WORD DERIVING DEVICE, KEY WORD DERIVING METHOD, AND  
STORAGE MEDIUM CONTAINING KEY WORD DERIVING PROGRAM



L E T T E R

Assistant Commissioner for Patents  
Box Patent Application  
Washington, D.C. 20231

October 19, 1999  
0397-0393P-SP

Sir:

Under the provisions of 35 USC 119 and 37 CFR 1.55(a), the applicant hereby claims the right of priority based on the following application(s):

<u>Country</u>	<u>Application No.</u>	<u>Filed</u>
JAPAN	10-300720	10/22/98

A certified copy of the above-noted application(s) is(are) attached hereto.

If necessary, the Commissioner is hereby authorized in this, concurrent, and future replies, to charge payment or credit any overpayment to deposit Account No. 02-2448 for any additional fees required under 37 C.F.R. 1.16 or under 37 C.F.R. 1.17; particularly, extension of time fees.

Respectfully submitted,

BIRCH, STEWART, KOLASCH & BIRCH, LLP

By:

TERRELL C. BIRCH

Reg. No. 19,382

P. O. Box 747

Falls Church, Virginia 22040-0747

Attachment  
(703) 205-8000  
/dl1

B.S.K.B.  
(703)205-8000  
NISHIMURA, Hideki  
390-393P  
10f1

# 日 本 国 特 許 庁

PATENT OFFICE  
JAPANESE GOVERNMENT

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日

Date of Application:

1998年10月22日

出 願 番 号

Application Number:

平成10年特許願第300720号

出 願 人

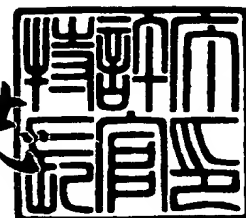
Applicant (s):

シャープ株式会社

1999年 8月 2日

特許庁長官  
Commissioner,  
Patent Office

伴佐山 建志



出証番号 出証特平11-3053924

【書類名】 特許願

【整理番号】 98-02657

【提出日】 平成10年10月22日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 17/30

【発明の名称】 キーワード抽出方法、キーワード抽出装置、及びキーワード抽出プログラムを記録したコンピュータ読み取り可能な記録媒体

【請求項の数】 7

【発明者】

    【住所又は居所】 大阪府大阪市阿倍野区長池町 2 2 番 2 2 号 シャープ株式会社内

    【氏名】 西村 英樹

【特許出願人】

    【識別番号】 000005049

    【氏名又は名称】 シャープ株式会社

    【電話番号】 06-621-1221

【代理人】

    【識別番号】 100103296

    【弁理士】

    【氏名又は名称】 小池 隆彌

    【電話番号】 06-621-1221

    【連絡先】 電話 0 4 3 - 2 9 9 - 8 4 6 6 知的財産権センター  
東京知的財産権部

【手数料の表示】

    【予納台帳番号】 012313

    【納付金額】 21,000円

【提出物件の目録】

    【物件名】 明細書 1

特平 10-300720

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9703283

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 キーワード抽出方法、キーワード抽出装置、及びキーワード抽出プログラムを記録したコンピュータ読み取り可能な記録媒体

【特許請求の範囲】

【請求項1】 データからキーワードを抽出する方法であって、  
所望のパラメータを用いて前記データを分割するステップと、  
分割されたグループ毎に単語を統計処理するステップと、  
統計処理された結果を比較し、重要度を算出するステップと、  
算出された重要度から、比較的重要度が高いと判断された単語からキーワード  
を決定するステップと、  
を備えたことを特徴とするキーワード抽出方法。

【請求項2】 前記所望のパラメータは、前記データにそれぞれ付加される  
属性から選択されてなることを特徴とする請求項1に記載のキーワード抽出方法  
。

【請求項3】 前記分割されたグループ毎に単語を統計処理する際、全ての  
グループに同一の統計処理を行ってなることを特徴とする請求項1又は2に記載  
のキーワード抽出方法。

【請求項4】 前記キーワードは、グループ毎のキーワード、又は全データ  
のキーワードとして決定されることを特徴とする請求項1～3のいずれかに記載  
のキーワード抽出方法。

【請求項5】 請求項4で決定されるキーワードは、パラメータとして選択  
された属性と関連づけて、これを全データの特性として判断されてなることを特  
徴とするキーワード抽出方法。

【請求項6】 データからキーワードを抽出する装置であって、  
所望のパラメータを用いて前記データを分割する手段と、  
分割されたグループ毎に単語を統計処理する手段と、  
統計処理された結果を比較し、重要度を算出する手段と、  
算出された重要度から、比較的重要度が高いと判断された単語からキーワード  
を決定する手段と、

を備えたことを特徴とするキーワード抽出装置。

【請求項 7】 コンピュータを、

所望のパラメータを用いて前記データを分割する手段、

分割されたグループ毎に単語を統計処理する手段、

統計処理された結果を比較し、重要度を算出する手段、

算出された重要度から、比較的重要度が高いと判断された単語からキーワードを決定する手段、

として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は大量のデータから、そのデータの特徴を示すキーワードを抽出する方法に関するもので、大量のデータの部分領域の統計処理を行うことによってキーワードを抽出するキーワード抽出方法およびキーワード検索装置およびキーワード抽出プログラムを記録した記録媒体に関する。

【0002】

【従来の技術】

大量のデータの部分領域の統計処理を行うことによってキーワードを抽出する方法が特開平 8-202737 号公報に記載されている。ここでは特許の明細書を例にあげ、まず、あらかじめ準備した、《発明の名称》、《特許請求の範囲》、等の見出し語に注目して全データを個々の段落に分割し、次に、同一文（センテンス）内での各単語の他の単語との共起数、同一段落毎の各単語の他の単語との共起数、全データでの各単語の出現数を求め、最後にこれらに適当な係数を乗じた代数和で各単語の重要度を求め、キーワードを決定している。

【0003】

すなわち、単なる各単語の出現頻度でキーワードを決定するのでは無く、同一文、同一段落で相互に共起する単語はより重要度（キーワードとしての妥当性）が高いと判断している。

【0004】

【発明が解決しようとする課題】

しかしながら、特開平 8-202737 号公報に記載の方法では、対象データの特殊性に基づいて予め準備された見出し語（《発明の名称》等）によって段落分割がなされるため、段落の分割が固定であった。また、抽出されるキーワードは対象データ全体に対するキーワードであって、個々の段落のキーワードは抽出されていなかった。

【0005】

従って、対象データが、特許明細書のように各段落がそれぞれ固定の意味を持ち、明細書 1 文書で内容が完結しているような場合には問題が少なかったが、対象データが各個人が送受信した（電子）メール全体とか、1 日、1 月単位のニュース全体とかのように、送受信相手、発生時刻（日時）、等の種々のパラメータで分割可能なデータ集合であり、対象データ全体の内容が把握しにくい場合には適用できなかった。

【0006】

本発明は、上記課題に基づいて創案されたもので、種々のパラメータで分割可能な大量のデータに対して各分割毎の統計処理結果の違いを比較することにより、各分割毎または全データのキーワードを抽出し、全データの特徴と、全データの中で各分割の特異性、傾向との両方または一方を把握することを目的とする。

【0007】

【課題を解決するための手段】

この発明（請求項 1）に係るキーワード抽出方法は、データからキーワードを抽出する方法であって、

所望のパラメータを用いて前記データを分割するステップと、

分割されたグループ毎に単語を統計処理するステップと、

統計処理された結果を比較し、重要度を算出するステップと、

算出された重要度から、比較的重要度が高いと判断された単語からキーワードを決定するステップと、

を備えたことにより、上記の目的を達成する。

【0008】

この発明（請求項2）に係るキーワード抽出方法は、請求項1において、  
前記所望のパラメータが、前記データにそれぞれ付加される属性から選択されてなることにより、上記の目的を達成する。

【0009】

この発明（請求項3）に係るキーワード抽出方法は、請求項1又は2において  
前記分割されたグループ毎に単語を統計処理する際、全てのグループに同一の統計処理を行ってなることにより、上記の目的を達成する。

【0010】

この発明（請求項4）に係るキーワード抽出方法は、請求項1～3のいずれかにおいて、  
前記キーワードが、グループ毎のキーワード、又は全データのキーワードとして決定されることにより、上記の目的を達成する。

【0011】

この発明（請求項5）に係るキーワード抽出方法は、  
請求項4で決定されるキーワードが、パラメータとして選択された属性と関連づけて、これを全データの特性として判断されてなることにより、上記の目的を達成する。

【0012】

この発明（請求項6）に係るキーワード抽出装置は、  
データからキーワードを抽出する装置であって、  
所望のパラメータを用いて前記データを分割する手段と、  
分割されたグループ毎に単語を統計処理する手段と、  
統計処理された結果を比較し、重要度を算出する手段と、  
算出された重要度から、比較的重要度が高いと判断された単語からキーワードを決定する手段と、  
を備えたことにより、上記の目的を達成する。



【0013】

この発明（請求項7）に係るキーワード抽出プログラムを記録したコンピュータ読み取り可能な記録媒体は、コンピュータを、  
所望のパラメータを用いて前記データを分割する手段、  
分割されたグループ毎に単語を統計処理する手段、  
統計処理された結果を比較し、重要度を算出する手段、  
算出された重要度から、比較的重要度が高いと判断された単語からキーワードを決定する手段、  
として機能させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体であることにより、上記の目的を達成する。

【0014】

即ち、本発明は、大量のデータを所望のパラメータを用いて、いくつかの数に分割するものである。この所望のパラメータは、データを社外向けメールと社内向けメールに分割する、メールの発生日または時刻で分割する、又は、特定の相手とそれ以外に分割する等、対象とするデータに付加された属性であれば、何でも良い。

【0015】

次に、分割されたグループ毎に単語の統計処理を行う。この統計処理は単なる単語の出現頻度でも良いし、他の単語との共起数であっても良い。

【0016】

更に、分割されたグループ毎の統計処理結果の比較を行う。比較を行う場合には、前記分割されたグループ毎の統計処理を、単なる出現頻度であれ、他の単語との共起数であれ、統一しておくことが好ましい。また、具体的な比較は統計結果の差分もしくは比率を求めることによってなされる。

【0017】

つまり、ある分割領域において、出現頻度又は他の単語との共起数が大きく、他の分割領域において、出現頻度又は他の単語との共起数が小さい単語は、重要度（キーワードとしての妥当性）が高いと判断する。

【0018】

最後に各分割で重要度が高いと判断された単語からキーワードを決定する。これは分割されたグループ毎のキーワードと判断しても良いし、全データのキーワードと判断しても良い。最も望ましいのは、最初に分割した時の属性と関連づけて、ある属性分割にはこのようなキーワードがあり、別の属性分割には別のキーワードがあるのが、この全データの特性であると判断することである。

【0019】

【発明の実施の形態】

以下、本発明の一実施例を添付図面に基づいて詳細に説明する。

【0020】

図1は本実施例に係るキーワード抽出装置の機能ブロック図である。

【0021】

この図において、1は対象データ選択手段である。対象データ選択手段1は、対象の文書集合に対して、各文書の単語やその他付加情報を整理する。2は対象データ分割手段である。対象データ分割手段2は、各文書の付加情報を利用して、文書集合を所望のN個に分割する。3は部分統計処理手段である。部分統計処理手段3は、分割されたN個の集合それぞれに対して単語の統計をとる。4は部分統計処理結果比較手段である。部分統計処理結果比較手段4は、N個の統計を元に、統計結果を比較し、違いを検出して、単語の重要度を決定する。5はキーワード抽出手段である。キーワード抽出手段5は、単語の重要度を参照し、ソートすることで、キーワードを抽出する。

【0022】

以後、ある開発資料の文書集合を対象に、著者を分割基準とした場合の動作を詳細に説明する。ここで、文書はコンピュータ上のファイルとなっており、著者情報を持ち、文書集合をファイル名のリストという形式で扱えるものとする。図2は、対象データ選択手段1における動作のフローチャートである。まず、文書集合から文書を1つ取り出す(S101)。取り出した文書を解析し、単語の出現頻度を調べる(S102)。この解析方法については、公知の技術を用いることにし、特に制限しない。

【0023】

次に、不要語辞書を用いて「そして、しかし、この」などの不要語を除去した後（S103）、著者名とともに文書テーブルに、文書の要素として追加する（S104）。文書集合の全ての文書に対して行うまで、S101～S104を繰り返す（S105）。結果として、文書テーブル6（図3）を得る。

【0024】

ここで、単語の後の（）の数字は、その単語の出現数を表す。なお、著者名は、ファイルの著者情報を利用するなど、特定できるものとする。ここでは、説明を簡単にするために各文書とも単語数が極端に少ない例にしているが、実際の文書では当然のことながら単語の種類・出現頻度とも多種・多数になる。

【0025】

図4は、対象データ分割手段2における動作のフローチャートである。文書テーブル6の全ての要素に対して、S201～S203を繰り返す。まず、文書テーブル6から、文書要素を1つ取り出す（S201）。例えば、「1 著者：幹部A 単語：画期的（10）、技術革新（5）、デジタル（4）」が取り出される。次に、著者名から、著者クラスを決定する（S202）。著者クラスとは、分割単位グループのことであり、クラスはいくつあっても良いが、この例では、「幹部」「技術」「企画」の3つのクラスがあるとする。従って、「1 著者：幹部A 単語：画期的（10）、技術革新（5）、デジタル（4）」の場合、クラス「幹部」と決定する。

【0026】

そして、決定されたクラスに従って、該当するクラスの単語リストに単語を追加する（S203）。「1 著者：幹部A 単語：画期的（10）、技術革新（5）、デジタル（4）」の場合、幹部クラスの単語リスト7に、単語「画期的（10）、技術革新（5）、デジタル（4）」を追加する（図5）。文書テーブル6の全ての要素に対して行うまで、S201～S203を繰り返す（S204）。

【0027】

結果として、「幹部」「企画」「技術」の3つのクラスに対応して、単語リス

ト7（図5）、単語リスト8（図6）、単語リスト9（図7）を得る。なお、著者クラスは著者名から直接決定するとしたが、著者クラステーブル10（図8）を用いて、クラスを決定するようにしても良い。

【0028】

図9は、部分統計処理手段3における動作のフローチャートである。全てのクラスに対して、S301～S306を繰り返す。まず、対象データ分割手段2によって得られた、あるクラスに対する、単語リストを選択する（S301）。例えば、幹部クラスの単語リスト7（図5）を選択する。選択した単語リストの全ての単語に対して、S302～S305を繰り返す。選択した単語リストから、単語を1つ取り出す（S302）。

【0029】

例えば、幹部クラスの単語リスト7から、単語「画期的（10）」を取り出す。そして、取り出した単語に対して、単語カウントテーブルに登録されているかどうかを調べ（S303）、登録されていれば、その単語に対応するカウンタを出現数だけ増加する（S304）。登録されていなければ、その単語を単語カウントテーブルに登録し、カウンタを出現数に設定する（S305）。

【0030】

例えば、幹部クラスを選択している時、単語リスト7から取り出された「画期的（10）」は1回目は単語カウントテーブル11に登録されていないので、単語カウントテーブル11に登録し、カウンタをその出現数である10にする。しかし、2回目に「画期的（2）」が取り出された時、単語カウントテーブル11に登録されているので、対応するカウンタを2増加し、12とする。

【0031】

S302～S305を、選択した単語リストの全ての単語に対して行うまで繰り返す（S306）。S301～S306を、全てのクラスに対して行うまで繰り返す（S307）。

【0032】

結果として、各クラスに対応した、単語カウントテーブル11、12、13を得る（図10、図11、図12）。

## 【0033】

図13は、部分統計処理結果比較手段4における動作のフローチャートである。S401～S404を、全てのクラスに対して行うまで繰り返す。まず、対象のクラスを選択する（S401）。例えば、幹部の観点で単語を抽出したい時、幹部クラスを選択する。S402～S403を、該当する単語カウントテーブルの全単語に対して行うまで、繰り返す。選択したクラスの単語カウントテーブルから単語を1つ選択する（S402）。例えば、技術クラスの単語カウントテーブル11から、単語「画期的」を選択する。

## 【0034】

次に、該当する単語カウントテーブルのカウンタ値と、他の単語カウントテーブルのカウンタ値を参照して、選択した単語の重要度を計算し、対応する単語重要度テーブルに登録する（S403）。この時、各クラスの文書数などを考慮に入れてもよい。

## 【0035】

例えば、単語kの重要度Iを以下の式によって求めることが出来る。

## 【0036】

【数1】

$$I(k) = \frac{\text{該当クラスにおける}k\text{のカウンタの値}}{\text{該当クラスの文書数}} \quad (1)$$

$$- \frac{\text{他のクラスにおける}k\text{のカウンタの値}}{\text{他のクラスの文書数}} \quad (2)$$

## 【0037】

このように重要度の決定に差分を用いると、カウンタ値の絶対値の大きさがI(k)の反映されることになり、より多く出現したキーワードを抽出したい場合に適している。S402～S403を、該当する単語カウントテーブルの全単語に対して行うまで、繰り返す（S404）。全てのクラスに対して、S401～S404を繰り返す（S405）。

## 【0038】

結果として、各クラスに対する単語重要度テーブル14、15、16を得る（図14、図15、図16）。

## 【0039】

図17は、キーワード抽出手段5における請求項1記載の方法のフローチャートである。S501～S503を、全てのキーワード抽出対象のクラスに対して行うまで繰り返す。まず、キーワード抽出対象のクラスを1つ選択する（S501）。そのクラスに対応する単語重要度テーブルを重要度でソートし、ソートした単語重要度テーブルを作成する（S502）。そして、ソートした単語重要度テーブル上位の単語を抽出キーワードとする（S503）。S501～S503を、全てのキーワード抽出対象のクラスに対して行うまで繰り返す（S504）。

## 【0040】

結果として、ソートした単語重要度テーブル17、18、19を得る（図18、図19、図20）。上述した、一連の処理を行うことで、対象の文書集合から各クラスの観点での特徴キーワードを抽出することが出来る。

## 【0041】

図21は、キーワード抽出手段5における請求項2記載の方法のフローチャートである。全てのクラスに対して、S511～S516を繰り返す。対象のクラスを1つ選択する（S511）。次に、選択クラスに対応する単語重要度テーブルの中の全ての単語に対して、S512～S517を繰り返す。対象の単語を1つ選択する（S512）。選択単語が単語統計テーブル20に登録されているかどうかを調べ（S513）、登録されていれば、単語統計テーブル20の対応する単語の統計値に、選択単語の重要度の絶対値を加え、さらに、選択クラスの単語カウントテーブルを参照し、対応する単語のカウンタに、選択単語のカウンタ値を加える（S514）。

## 【0042】

登録されていなければ、単語統計テーブル20に選択単語を登録した後（S515）、単語統計テーブル20の対応する単語の統計値に、選択単語の重要度の絶対値を設定し、さらに、選択クラスの単語カウントテーブルを参照し、対応する単語のカウンタに、選択単語のカウンタ値を設定する（S516）。統計値に対しては、重要度の自乗などを加えるようにしても良い。S512～S516を

、全ての単語に対して行うまで繰り返す（S517）。S511～S517を、全てのクラスに対して行うまで繰り返す（S518）。

#### 【0043】

最後に、単語統計テーブル20を統計値が小さく、カウンタ値が大きい順にソートし、ソートした単語統計テーブル21を作成する（S519）。この時、具体的には、例えば

係数C × 統計値 - カウンタ値

の値を基準にソート、統計値とカウンタ値のどちらかに閾値を設定することで絞り込み残りの値でソート、統計値／カウンタ値の値を基準にソート、などとすれば良い。そして、ソートした単語統計テーブル21の上位単語を、抽出キーワードとする（S520）。

#### 【0044】

結果として得られた、ソートした単語統計テーブル21を図22に示す。ここでは、統計値／カウンタ値の値の小さい順にソートした。上述した、一連の処理を行うことで、対象の文書集合から各クラスによって出現傾向の似ているキーワードを抽出することが出来る。

#### 【0045】

本実施例では、該当クラスと他のクラスとの差分に着目したが、比率に着目することも出来る。すなわち、図13における、S403の重要度Iの計算において、

#### 【0046】

【数2】

$$I(k) = \frac{\text{該当クラスにおける} k \text{のカウンタの値}}{\text{該当クラスの文書数}} \quad (3)$$

$$\div \frac{\text{他のクラスにおける} k \text{のカウンタの値}}{\text{他のクラスの文書数}} \quad (4)$$

#### 【0047】

とすればよい。このように重要度の決定に比率を用いると、文書数の絶対数に関わらず重要度の大小の幅が大きくなり、文書数によらずに重要度を考慮したキー

ワードを抽出したい場合に適している。

【0048】

この時、統計値の計算をする時に、無限大なる場合を考慮し、重要度に上限、下限を設定し、さらに、重要度が1未満ならその逆数を重要度の代わりに使うようにするか、対数の絶対値を使うなど、重要度の違いが累積されるようにする。または、重要度が1未満の時その逆数を重要度の代わりに使う場合には、重要度またはその逆数が1以上になることに注目し、統計値に加えるのではなく、乗じることにも出来る。

【0049】

また、本実施例では、対象の文書集合を開発資料とし、分割基準を著者としたが、その他分割基準になりうる所望の基準で良く、例えば、メールを対象として、社外用のメールと、社内用のメールに分割したり、メールの発生日または時刻で分割したり、特定の相手とそれ以外に分割するなど、対象データに付加されているものであれば何でも良い。

【0050】

様々な検索システムでは、利用者のアクセスをログとして記録している場合が多い。この時、アクセスログにはアクセスした利用者のIDや、アクセス日時や、検索キー（多くの場合キーワード）が記録されている。従って、この一回のアクセスに対するログを1文書として扱くと、アクセスログ全体で、文書集合を構成することになり、利用者IDで分割したり、アクセス日時を用いて昼夜あるいは平日、休日などに分割することで、それぞれのアクセス傾向を調べることが出来る。

【0051】

【発明の効果】

本発明のキーワード抽出方法、キーワード抽出装置、又は記録媒体を用いることによれば、種々のパラメータで分割可能な大量のデータに対して、分割されたグループ毎の統計処理結果の違いを比較することによって、分割されたグループ毎、又は全データのキーワードを抽出し、全データの特徴と、全データの中での分割されたグループの特異性及び／又は傾向を把握することが出来、しかも、デ



一タの形式として特定のものは必要でなく、所望のものを扱える。

【0052】

特に、重要度の決定に、出現頻度の差分を用いる方法では、より出現頻度の高いキーワードを抽出することが出来、重要度の決定に、出現頻度の比率を用いる方法では、文書の絶対数によらずに重要度を考慮したキーワードを抽出することが出来る。

【図面の簡単な説明】

【図1】

本発明の機能ブロック図である。

【図2】

データ対象選択手段における動作のフローチャートである。

【図3】

文書テーブルの例を示す説明図である。

【図4】

対象データ分割手段における動作のフローチャートである。

【図5】

幹部クラスの単語リストの例を示す説明図である。

【図6】

企画クラスの単語リストの例を示す説明図である。

【図7】

技術クラスの単語リストの例を示す説明図である。

【図8】

著者クラステーブルの例を示す説明図である。

【図9】

部分統計処理手段における動作のフローチャートである。

【図10】

幹部クラスの単語カウントテーブルの例を示す説明図である。

【図11】

企画クラスの単語カウントテーブルの例を示す説明図である。

【図 1 2】

技術クラスの単語カウントテーブルの例を示す説明図である。

【図 1 3】

部分統計処理結果比較手段における動作のフローチャートである。

【図 1 4】

幹部クラスの単語重要度テーブルの例を示す説明図である。

【図 1 5】

企画クラスの単語重要度テーブルの例を示す説明図である。

【図 1 6】

技術クラスの単語重要度テーブルの例を示す説明図である。

【図 1 7】

キーワード抽出手段における請求項 1 記載の方法のフローチャートである。

【図 1 8】

幹部クラスのソートした単語重要度テーブルの例を示す説明図である。

【図 1 9】

企画クラスのソートした単語重要度テーブルの例を示す説明図である。

【図 2 0】

技術クラスのソートした単語重要度テーブルの例を示す説明図である。

【図 2 1】

キーワード抽出手段における請求項 2 記載の方法のフローチャートである。

【図 2 2】

単語統計テーブルの例を示す説明図である。

【図 2 3】

ソートした単語統計テーブルの例を示す説明図である。

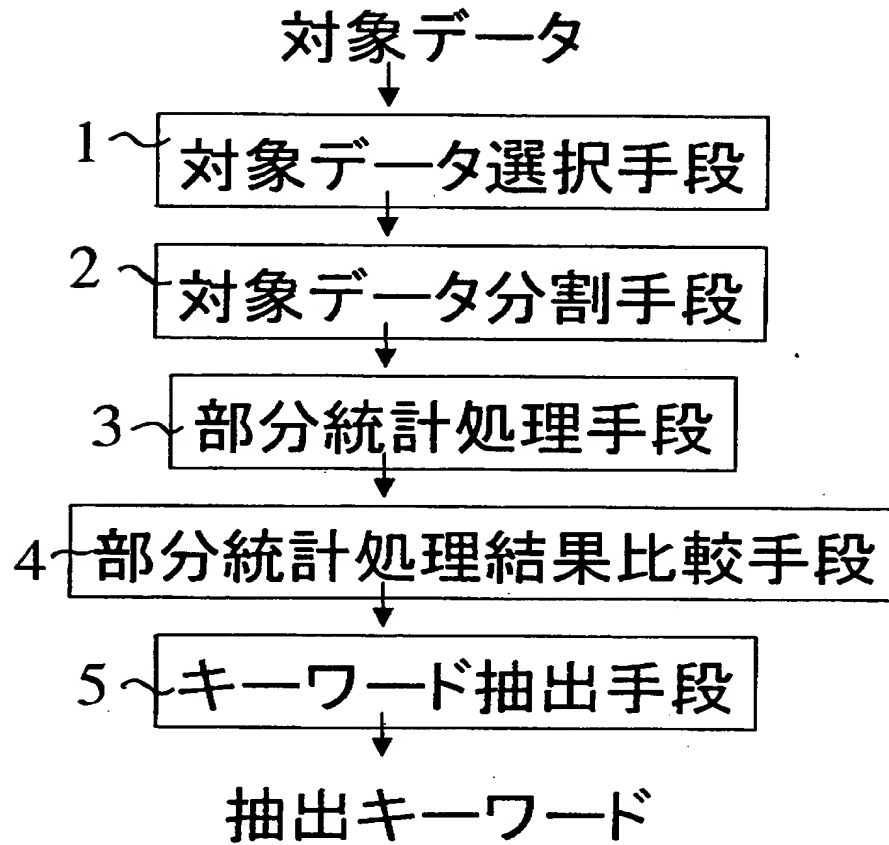
【符号の説明】

- 1 対象データ選択手段
- 2 対象データ分割手段
- 3 部分統計処理手段
- 4 部分統計処理結果比較手段

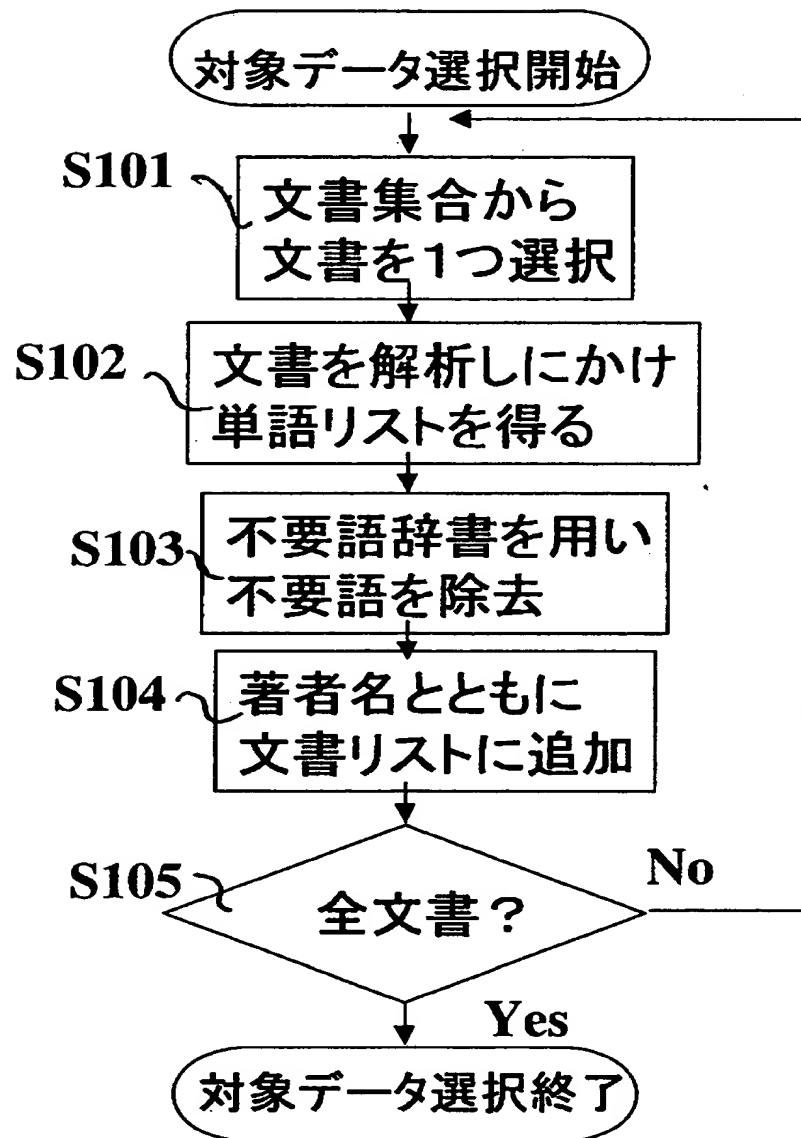
- 5 キーワード抽出手段
- 6 文書テーブル
- 7 幹部クラスの単語リスト
- 8 企画クラスの単語リスト
- 9 技術クラスの単語リスト
- 10 著者クラステーブル
- 11 幹部クラスの単語カウントテーブル
- 12 企画クラスの単語カウントテーブル
- 13 技術クラスの単語カウントテーブル
- 14 幹部クラスの単語重要度テーブル
- 15 企画クラスの単語重要度テーブル
- 16 技術クラスの単語重要度テーブル
- 17 幹部クラスのソートした単語重要度テーブル
- 18 企画クラスのソートした単語重要度テーブル
- 19 技術クラスのソートした単語重要度テーブル
- 20 単語統計テーブル
- 21 ソートした単語統計テーブル

【書類名】 図面

【図 1】



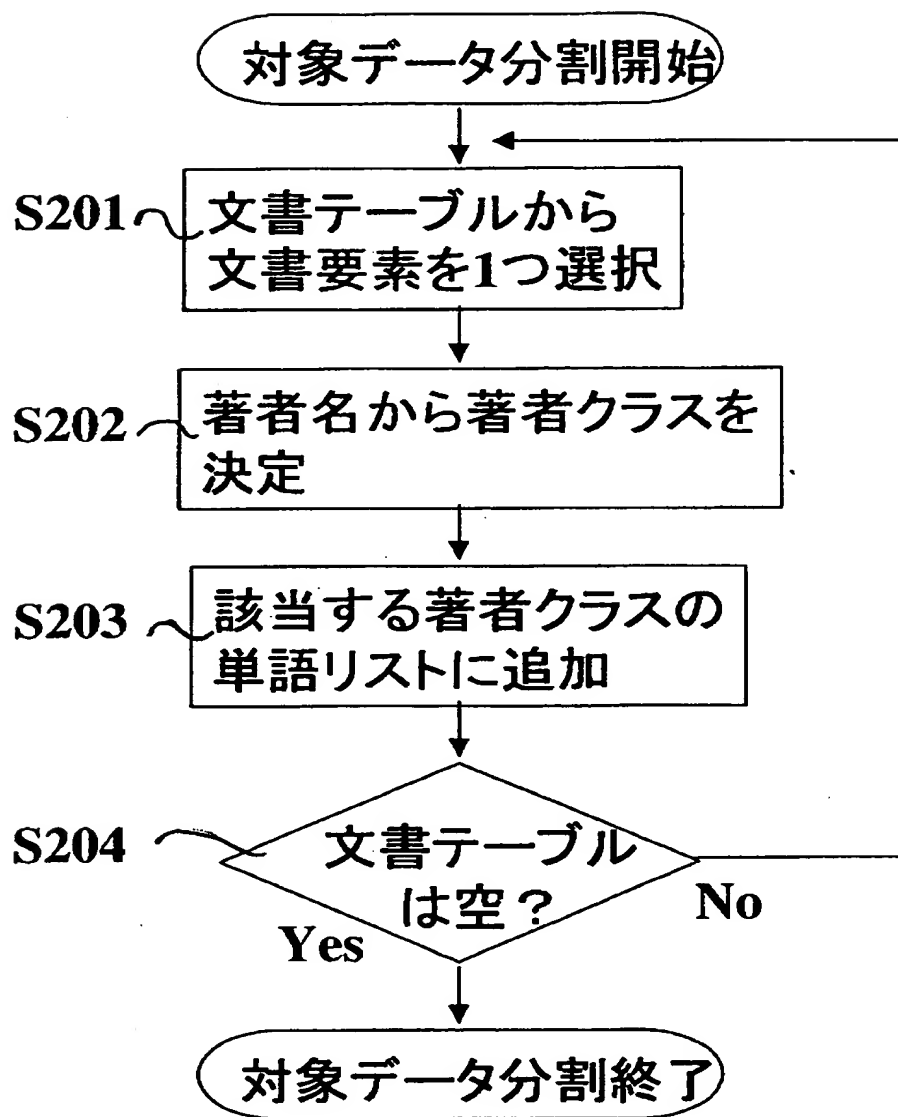
【図2】



【図 3】

文書番号	著者	単語 (出現数)
1	幹部 A	画期的 (10), 技術革新 (5), デジタル (4)
2	企画 A	多彩 (7), デジタル (3), 高画質 (2), 画期的 (2)
3	技術 A	デジタル (6), 通信 (5), 詳細仕様 (3)
4	技術 B	デジタル (4), 通信 (8)
5	幹部 B	画期的 (8), デジタル (6), 技術革新 (7), 情報化社会 (5), 高画質 (4)
6	技術 C	技術革新 (3), 高画質 (6)
7	幹部 B	技術革新 (6), 経営環境 (4), 情報化社会 (5)
8	企画 B	多彩 (3), 高画質 (3), 鮮やか (4), 通信 (2)
9	幹部 C	経営環境 (5), 通信 (2), デジタル (8)
10	幹部 C	画期的 (7), 技術革新 (4)
11	技術 C	詳細仕様 (3), デジタル (8), プロトコル (5)
12	企画 C	多彩 (6), デジタル (4), 鮮やか (4)

【図4】



【図 5】

単語
画期的 (10), 技術革新 (5), デジタル (4),
画期的 (8), デジタル (6), 技術革新 (7), 情報化社会 (5), 高画質 (4),
技術革新 (6), 経営環境 (4), 情報化社会 (5),
経営環境 (5), 通信 (2), デジタル (8),
画期的 (7), 技術革新 (4)

【図 6】

単語
多彩 (7), デジタル (3), 高画質 (2), 画期的 (2),
多彩 (3), 高画質 (3), 鮮やか (4), 通信 (2),
多彩 (6), デジタル (4), 鮮やか (4)

【図 7】

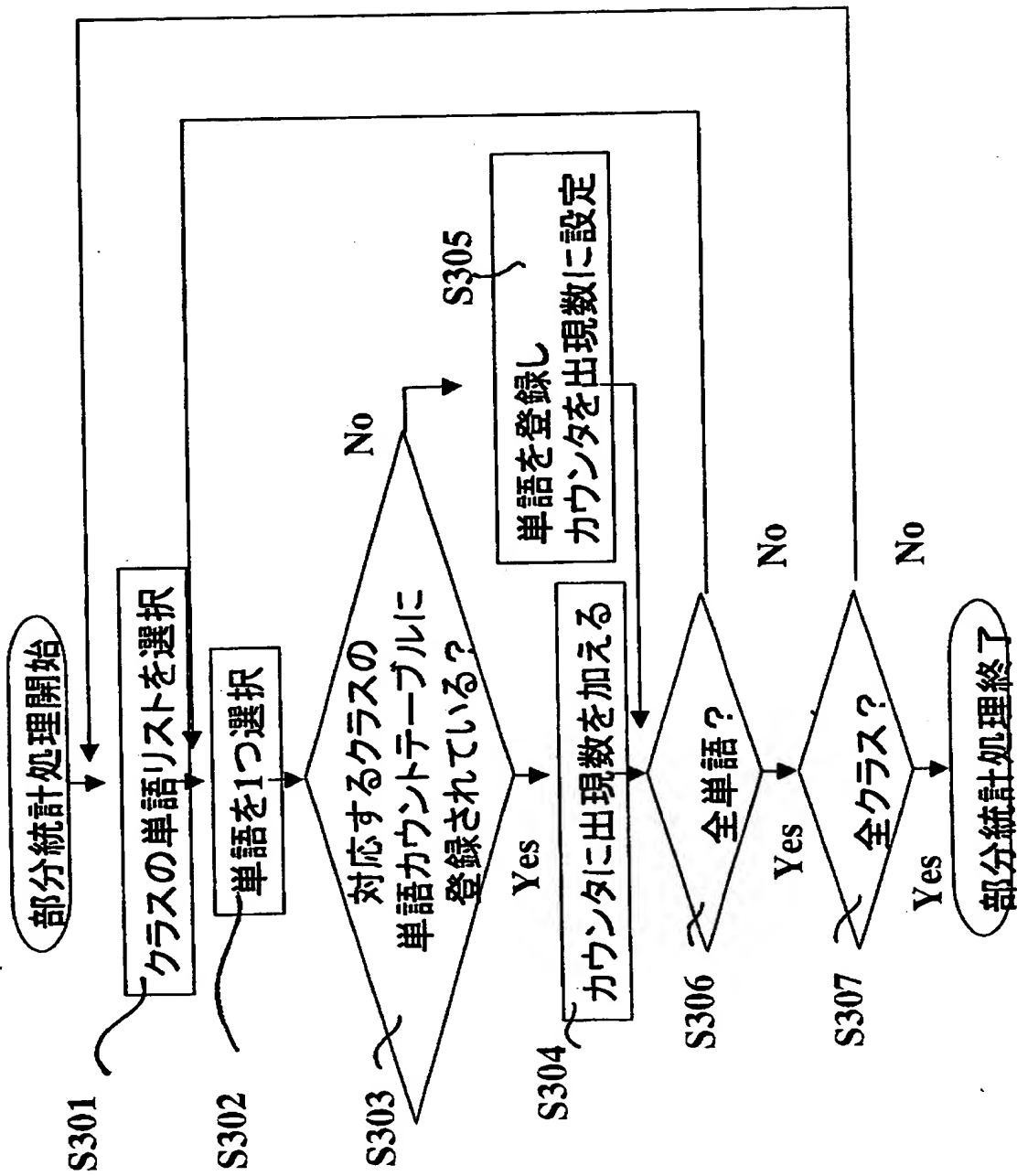
単語
デジタル (6), 通信 (5), 詳細仕様 (3),
デジタル (4), 通信 (8),
技術革新 (3), 高画質 (6),
詳細仕様 (3), デジタル (8), プロトコル (5)



【図 8】

著者名	著者クラス
幹部 A	幹部
幹部 B	幹部
幹部 C	幹部
企画 A	企画
企画 B	企画
企画 C	企画
技術 A	技術
技術 B	技術
技術 C	技術

【図9】



【図 10】

単語	カウンタ
画期的	25
技術革新	22
デジタル	18
情報化社会	10
経営環境	9
高画質	4
通信	2

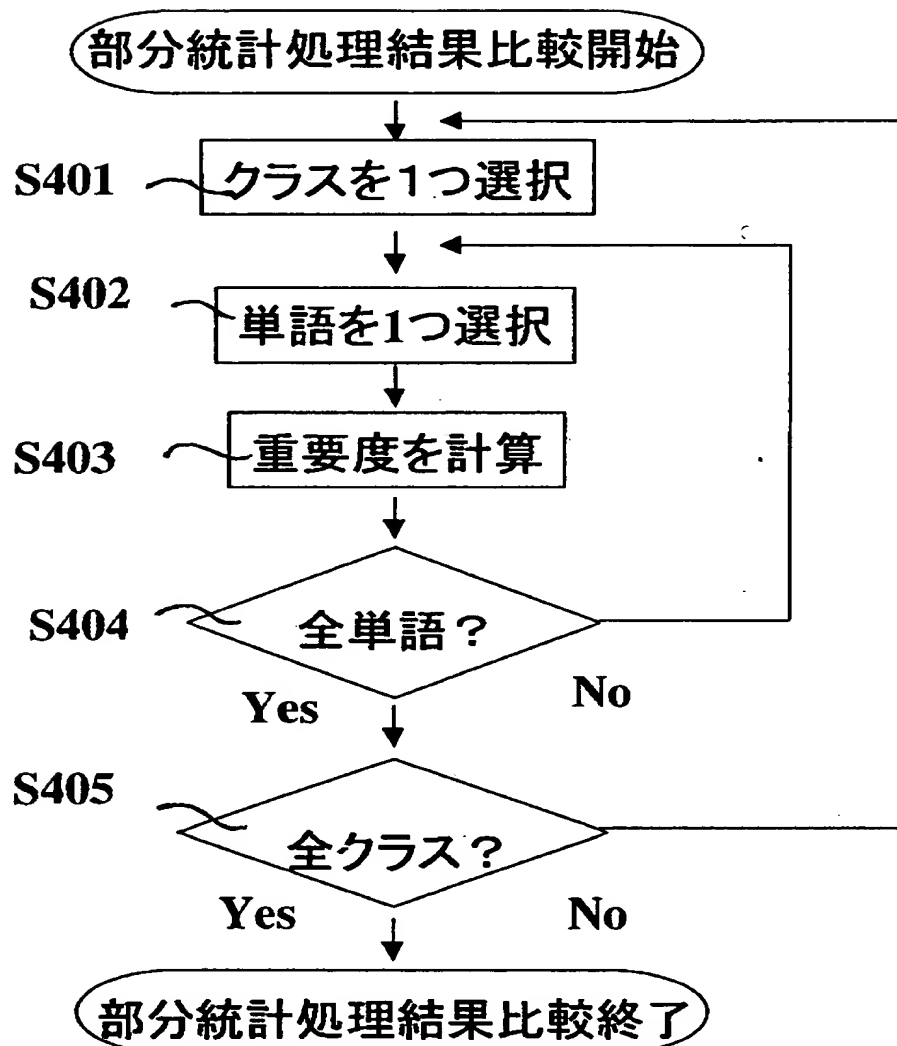
【図 11】

単語	カウンタ
多彩	16
鮮やか	8
デジタル	7
高画質	5
画期的	2
通信	2

【図 12】

単語	カウンタ
デジタル	18
通信	13
詳細仕様	6
技術革新	3
高画質	6
プロトコル	5

【図 13】



【図 14】

単語	重要度
画期的	4.714
技術革新	3.971
経営環境	1.800
デジタル	0.028
情報化社会	2.000
高画質	-0.771
通信	-1.742

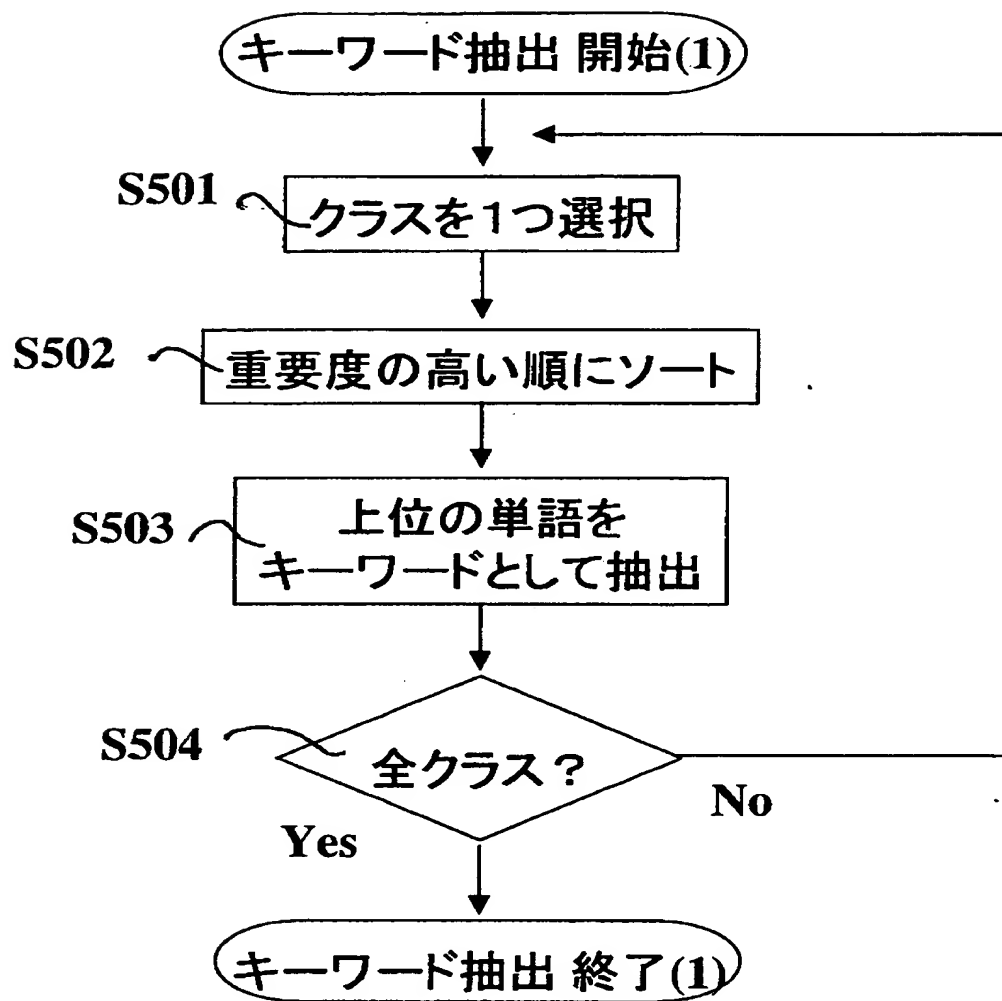
【図 15】

単語	重要度
多彩	5.333
高画質	0.555
鮮やか	2.666
デジタル	-1.666
画期的	-2.111
通信	-1.000

【図 16】

単語	重要度
デジタル	1.375
詳細仕様	1.500
通信	2.750
技術革新	-2.000
高画質	0.375
プロトコル	1.250

【図 17】



【図 18】

単語	重要度
画期的	4.714
技術革新	3.971
情報化社会	2.000
経営環境	1.800
デジタル	0.028
高画質	-0.771
通信	-1.742

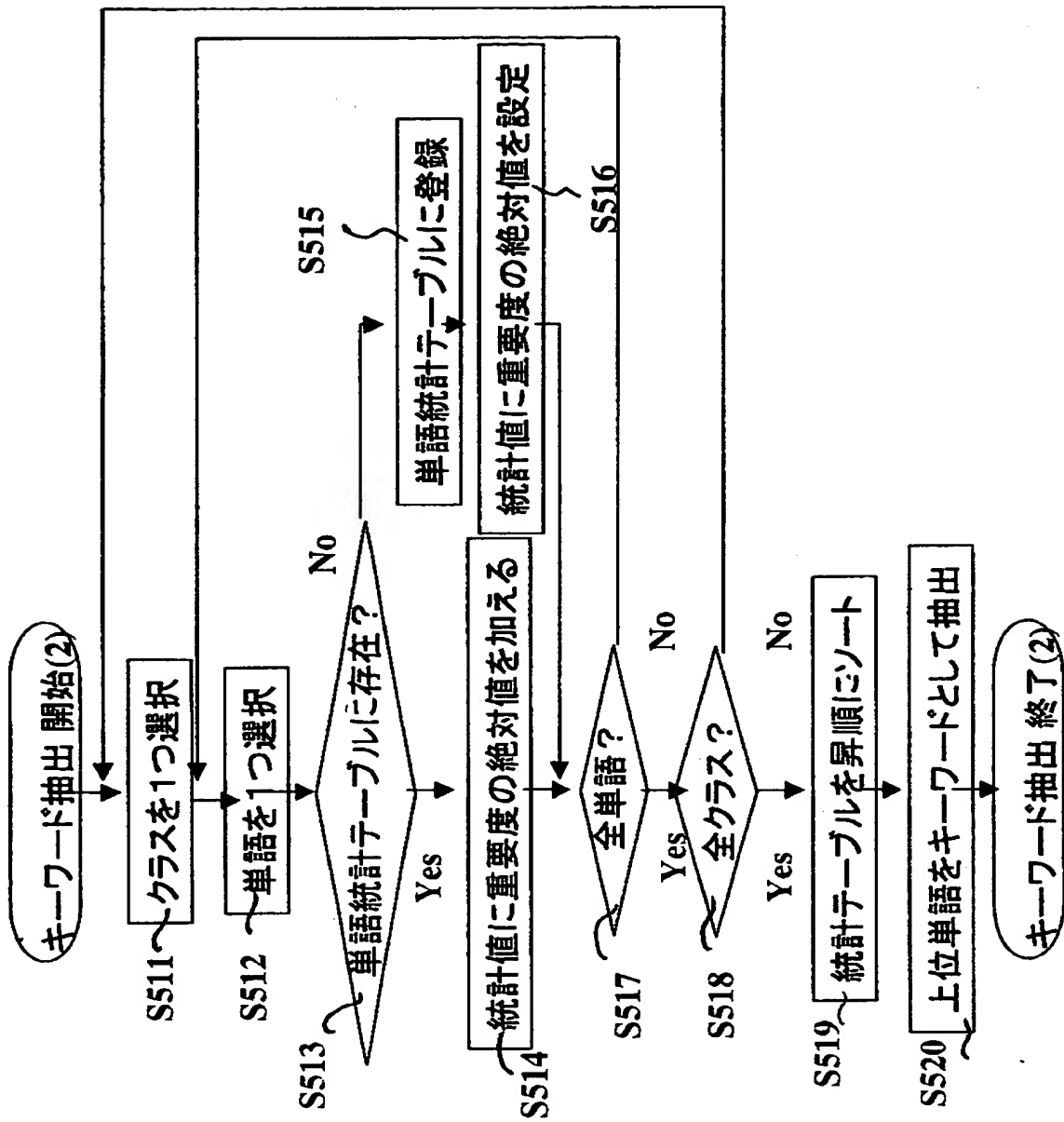
【図 19】

単語	重要度
多彩	5.333
鮮やか	2.666
高画質	0.555
通信	-1.000
デジタル	-1.666
画期的	-2.111

【図 20】

単語	重要度
通信	2.750
デジタル	1.375
詳細仕様	1.500
プロトコル	1.250
高画質	0.375
技術革新	-2.000

【図 21】





【図 22】

単語	統計値	カウンタ
画期的	6.825	27
技術革新	5.971	25
情報化社会	2.000	10
経営環境	1.800	9
デジタル	3.069	43
高画質	1.701	15
通信	5.492	17
詳細仕様	1.500	6
多彩	5.333	16
鮮やか	2.666	8
プロトコル	1.250	5

【図 23】

単語	統計値	カウンタ	統計値/カウンタ
デジタル	3.069	43	0.071
高画質	1.701	15	0.113
情報化社会	2.000	10	0.200
経営環境	1.800	9	0.200
技術革新	5.971	25	0.238
詳細仕様	1.500	6	0.250
プロトコル	1.250	5	0.250
画期的	6.825	27	0.252
通信	5.492	17	0.323
多彩	5.333	16	0.333
鮮やか	2.666	8	0.333

【書類名】 要約書

【要約】

【課題】 キーワードを抽出したいデータの形状に依存せず、データを所望のパラメータで分割し、キーワードを抽出する。

【解決手段】 データからキーワードを抽出する方法であって、所望のパラメータを用いて前記データを分割するステップと、分割されたグループ毎に単語を統計処理するステップと、統計処理された結果を比較し、重要度を算出するステップと、算出された重要度から、比較的重要度が高いと判断された単語からキーワードを決定するステップと、を備えたことを特徴とする。

【選択図】 なし

【書類名】 職権訂正データ  
【訂正書類】 特許願

<認定情報・付加情報>

【特許出願人】

【識別番号】 000005049

【住所又は居所】 大阪府大阪市阿倍野区長池町 2 2 番 2 2 号

【氏名又は名称】 シャープ株式会社

【代理人】 申請人

【識別番号】 100103296

【住所又は居所】 大阪府大阪市阿倍野区長池町 2 2 番 2 2 号 シャー  
プ株式会社内

【氏名又は名称】 小池 隆彌

特平 10-300720

出 願 人 履 歴 情 報

識別番号 [000005049]

1. 変更年月日	1990年 8月29日
[変更理由]	新規登録
住 所	大阪府大阪市阿倍野区長池町22番22号
氏 名	シャープ株式会社

BEST AVAILABLE COPY